

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2003-022088

(43)Date of publication of application : 24.01.2003

(51)Int.Cl.

G10L 15/06
G06F 3/16
G10L 15/02
G10L 15/10
G10L 15/14
G10L 15/18

(21)Application number : 2001-209503

(71)Applicant : SHARP CORP

(22)Date of filing : 10.07.2001

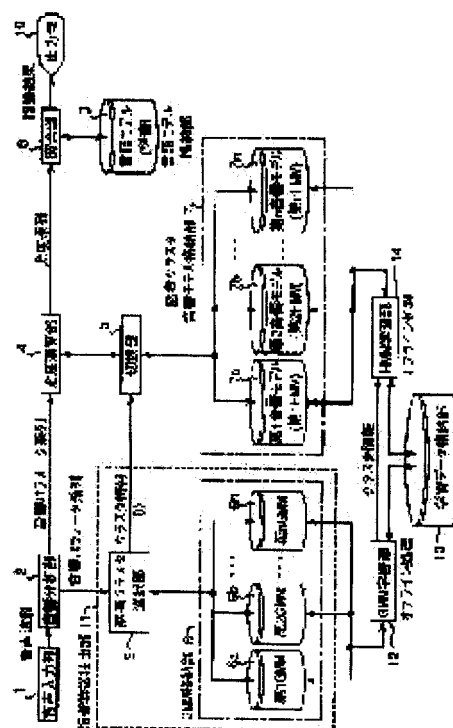
(72)Inventor : YAMAGUCHI KOICHI

(54) DEVICE AND METHOD FOR SPEAKER'S FEATURES EXTRACTION, VOICE RECOGNITION DEVICE, AND PROGRAM RECORDING MEDIUM

(57)Abstract:

PROBLEM TO BE SOLVED: To precisely extract speaker's features from a less amount of speech data.

SOLUTION: A GMM learning part 12 adds the value of a vocal tract length expansion/contraction coefficient α to voice data of respective learners stored in a learning data storage part 13, clusters the learners according to the vocal track expansion/contraction coefficient α , performs data conversion so that voice data of a speaker in a cluster nearby some cluster C belong to the cluster C, and reclusters the learners by using GMMs of the respective clusters. The GMMs of the obtained (n) clusters are stored in a GMM storage part 6. A speaker cluster selection part 3 makes the (n) GMMs stored in the GMM storage part 6 operate on a sound parameter series from a sound analysis part 2 and outputs the index of the GMM giving the maximum likelihood as speaker cluster information. Thus, speaker's features are precisely extracted from a small amount of learning data without depending upon speech contents.



LEGAL STATUS

[Date of request for examination]

18.06.2004

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(51)Int.Cl. ⁷	識別記号	F I	ページ* (参考)
G 1 0 L 15/06		G 0 6 F 3/16	3 2 0 G 5 D 0 1 5
G 0 6 F 3/16	3 2 0		3 2 0 H
		G 1 0 L 3/00	5 2 1 V
G 1 0 L 15/02			5 1 5 E
15/10			5 2 1 S
審査請求 未請求 請求項の数12 O L (全 16 頁) 最終頁に続く			

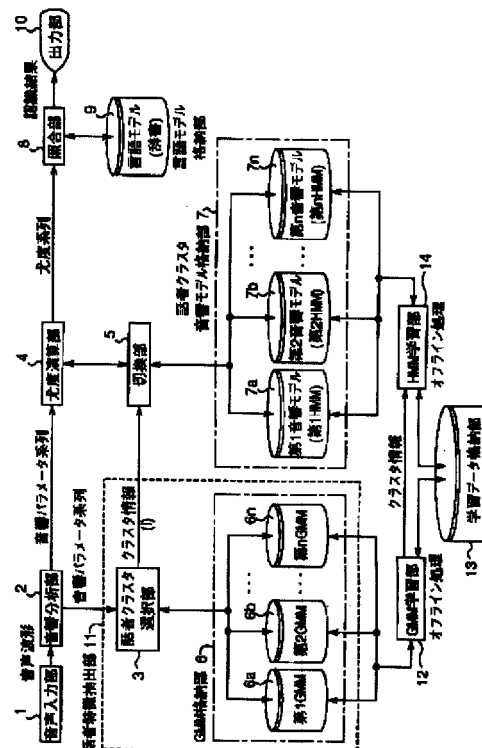
(21)出願番号	特願2001-209503(P2001-209503)	(71)出願人	000005049 シャープ株式会社 大阪府大阪市阿倍野区長池町22番22号
(22)出願日	平成13年7月10日(2001.7.10)	(72)発明者	山口 耕市 大阪府大阪市阿倍野区長池町22番22号 シャープ株式会社内
		(74)代理人	100062144 弁理士 青山 葆 (外1名)
		Fターム(参考)	5D015 FF04 GG04 HH04

(54) 【発明の名称】 話者特徴抽出装置および話者特徴抽出方法、音声認識装置、並びに、プログラム記録媒体

(57) 【要約】

【課題】 より少ない発声データから精度良く話者特徴を抽出する。

【解決手段】 GMM学習部12は、学習データ格納部13に格納された各学習話者の音声データに声道長伸縮係数 α の値を与え、この声道長伸縮係数 α に基づいて学習話者をクラスタリングし、あるクラスタCの近傍のクラスタDに属する話者の音声データを上記クラスタCに属するようにデータ変換し、各クラスタのGMMを用いて学習話者をクラスタリングし直す。得られたn個のクラスタのGMMはGMM格納部6に格納される。話者クラスタ選択部3は、音響分析部2からの音響パラメータ系列にGMM格納部6に格納されたn個のGMMを作用させ、最大尤度を与えるGMMのインデックスを話者クラスタ情報として出力する。こうして、少ない学習データから、発話内容に因らずに精度良く話者特徴を抽出する。



【特許請求の範囲】

【請求項1】 入力話者の音声に基づいて、標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として抽出する話者特徴抽出装置において、各学習話者に関して、上記標準話者に対する声道長の伸縮係数 α を所定の方法によって予め求め、この求められた伸縮係数 α の値に基づいて上記学習話者をクラスタリングする学習話者クラスタリング手段と、上記クラスタリングされた各クラスタに属する話者集合毎に、学習によって1状態の混合ガウス分布型音響モデルを生成する音響モデル生成手段と、上記生成された1状態の混合ガウス分布型音響モデルの夫々に対する上記学習話者の音声サンプルの尤度を算出し、その尤度に基づいて上記学習話者を再クラスタリングする再クラスタリング手段と、上記音響モデル生成手段と再クラスタリング手段とを制御して、所定の条件を満たすまで、上記1状態の混合ガウス分布型音響モデルの生成と上記学習話者の再クラスタリングとを繰り返すループ学習手段と、上記ループ学習手段によって最終的に生成された1状態の混合ガウス分布型音響モデルの群を格納する音響モデル格納部と、上記音響モデル格納部に格納された1状態の混合ガウス分布型音響モデルの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈する1状態の混合ガウス分布型音響モデルを入力話者が属するクラスタの情報として選択する話者クラスタ選択部を備えて、上記入力話者の特徴として上記最大の尤度を呈する1状態の混合ガウス分布型音響モデルを抽出することを特徴とする話者特徴抽出装置。

【請求項2】 請求項1に記載の話者特徴抽出装置において、上記学習話者クラスタリング手段によってクラスタリングされた各クラスタのうちの注目クラスタに隣接する隣接クラスタに属する学習話者の音声サンプル、または、上記注目クラスタと上記伸縮係数 α 値の差が所定値以内の近傍クラスタに属する学習話者の音声サンプルに対して、上記注目クラスタと隣接クラスタまたは近傍クラスタとの上記伸縮係数 α 値に基づいて周波数伸縮を行うことによって、上記注目クラスタに属する音声サンプルを生成し、この生成された音声サンプルを上記注目クラスタに編入して当該クラスタの音声サンプル数を豊富化する操作を、上記学習話者クラスタリング手段によってクラスタリングされた総てのクラスタについて実行する音声サンプル豊富化手段を備えて、上記音響モデル生成手段は、上記音声サンプル豊富化手段によって音声サンプル数が豊富化された後の各クラスタ毎に、上記1状態の混合ガウス分布型音響モデルを生成するようになっていることを特徴とする話者特徴抽出装置。

【請求項3】 請求項2に記載の話者特徴抽出装置において、上記注目クラスタに属する音声サンプルを生成する際に、上記音声サンプル豊富化手段が上記学習話者の音声サンプルに対して周波数伸縮を行う音声区間を、有音・無音の別および調音点に基づいて限定するようにしたことを特徴とする話者特徴抽出装置。

【請求項4】 請求項1あるいは請求項2に記載の話者特徴抽出装置において、上記再クラスタリング手段によって上記学習話者を再クラスタリングする場合に、上記再クラスタリングの対象となる対象学習話者が再クラスタリングの前に属していたクラスタの伸縮係数 α と再クラスタリング後に属するクラスタの伸縮係数 α とが所定値以上離れている場合には、当該対象学習話者を上記再クラスタリングの対象から外すようになっていることを特徴とする話者特徴抽出装置。

【請求項5】 請求項1あるいは請求項2に記載の話者特徴抽出装置において、上記ループ学習手段によって最終的にクラスタリングされた各クラスタに属する学習話者を更にクラスタリングしてサブクラスタを生成し、上記各サブクラスタに属する話者集合毎に学習によって1状態の混合ガウス分布型音響モデルを生成するサブクラスタ生成手段を備えて、上記音響モデル格納部は、上記サブクラスタ生成手段によって生成された1状態の混合ガウス分布型音響モデルの群を、各サブクラスタの伸縮係数 α に対応付けて格納するようになっていることを特徴とする話者特徴抽出装置。

【請求項6】 入力話者の音声に基づいて、標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として抽出する話者特徴抽出装置において、上記標準話者に対する声道長の伸縮係数 α の値に基づいて学習話者をクラスタリングし、各クラスタに属する話者集合毎に1状態の混合ガウス分布型音響モデルを生成し、この生成された1状態の混合ガウス分布型音響モデルの夫々に対する上記学習話者の音声サンプルの尤度に基づいて上記学習話者を再クラスタリングし、上記1状態の混合ガウス分布型音響モデルの生成と上記学習話者の再クラスタリングとを所定の条件を満たすまで繰り返して最終的に生成された1状態の混合ガウス分布型音響モデルの群が格納された音響モデル格納部と、上記音響モデル格納部に格納された1状態の混合ガウス分布型音響モデルの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈する1状態の混合ガウス分布型音響モデルを入力話者が属するクラスタの情報として選択する話者クラスタ選択部を備えて、上記入力話者の特徴として上記最大の尤度を呈する1状態の混合ガウス分布型音響モデルを抽出することを特徴とする話者特徴抽出装置。

【請求項7】 音響モデルとして隠れマルコフモデルを用い、入力話者の音声に基づいて抽出された標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として、上記入力話者の音声を認識する音声認識装置であって、

請求項1あるいは請求項6に記載の話者特徴抽出装置と、

上記話者特徴抽出装置の音響モデル格納部に格納された各1状態の混合ガウス分布型音響モデルによって表わされる話者クラスタに属する話者集合毎に、学習によって生成された隠れマルコフモデルの群を格納する隠れマルコフモデル格納部と、

上記話者特徴抽出装置によって選択されたクラスタに基づいて、上記隠れマルコフモデル格納部に格納されている上記選択されたクラスタに対応するクラスタの隠れマルコフモデルを、音声認識用の音響モデルとして切り換え選出する切換部を備えたことを特徴とする音声認識装置。

【請求項8】 音響モデルとして隠れマルコフモデルを用い、入力話者の音声に基づいて抽出された標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として、上記入力話者の音声を認識する音声認識装置であって、

請求項1あるいは請求項6に記載の音響モデル格納部に格納された各1状態の混合ガウス分布型音響モデルによって表わされる話者クラスタに属する話者集合毎に、学習によって生成された隠れマルコフモデルの群を格納する隠れマルコフモデル格納部と、

上記隠れマルコフモデル格納部に格納された隠れマルコフモデルの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈する隠れマルコフモデルを音声認識用の音響モデルとして切り換え選出する切換部を備えたことを特徴とする音声認識装置。

【請求項9】 声道長の伸縮関数 α を用いて入力音声のスペクトルの周波数軸を伸縮することによって入力話者の音響特徴量を標準話者の音響特徴量に正規化する話者正規化手段を有する音声認識装置において、

上記話者正規化手段は、

請求項1あるいは請求項6に記載の話者特徴抽出装置と、

上記入力話者の音声サンプルに基づいて、上記話者特徴抽出装置によって選択された1状態の混合ガウス分布型音響モデルに対応する声道長の伸縮係数 α を用いて、上記入力音声のスペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

【請求項10】 声道長の伸縮関数 α を用いて音声のスペクトルの周波数軸を伸縮することによって、音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、

上記話者適応手段は、

請求項1あるいは請求項6に記載の話者特徴抽出装置と、

上記入力話者の音声サンプルに基づいて、上記話者特徴抽出装置によって選択された1状態の混合ガウス分布型音響モデルに対応する声道長の伸縮係数 α の逆数を用いて、上記音響モデルのスペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

10 【請求項11】 入力話者の音声に基づいて、標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として抽出する話者特徴抽出方法であって、

各学習話者に関して、上記標準話者に対する声道長の伸縮係数 α を所定の方法によって求め、この求められた伸縮係数 α の値に基づいて上記学習話者をクラスタリングし、

上記クラスタリングされた各クラスタに属する話者集合毎に、学習によって1状態の混合ガウス分布型音響モデルを生成し、

20 上記生成された1状態の混合ガウス分布型音響モデルの夫々に対する上記学習話者の音声サンプルの尤度を算出し、その尤度に基づいて上記学習話者を再クラスタリングし、

所定の条件を満たすまで、上記1状態の混合ガウス分布型音響モデルの生成と上記学習話者の再クラスタリングとを繰り返すループ学習を行い、

30 上記ループ学習によって最終的に生成された1状態の混合ガウス分布型音響モデルの群を音響モデル格納部に格納し、

上記音響モデル格納部に格納された1状態の混合ガウス分布型音響モデルの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈する1状態の混合ガウス分布型音響モデルを上記入力話者の特徴として抽出することを特徴とする話者特徴抽出方法。

【請求項12】 コンピュータを、

40 請求項1に記載の学習話者クラスタリング手段、音響モデル生成手段、再クラスタリング手段、ループ学習手段、音響モデル格納部および話者クラスタ選択部として機能させる話者特徴抽出処理プログラムが記録されたことを特徴とするコンピュータ読出し可能なプログラム記録媒体。

【発明の詳細な説明】

【0001】

50 【発明の属する技術分野】この発明は、標準話者の音声スペクトルに対する入力音声スペクトルの周波数軸の線形伸縮係数を話者特徴として抽出する話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、音声合成装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体に関する。

【0002】

【従来の技術】従来より、隠れマルコフモデル(Hidden Markov Model: 以下、HMMと言う)を用いた音声認識方法の開発が盛んに行われている。このHMMは、大量の音声データから得られる音声の統計的特徴を確率的にモデル化したものであり、このHMMを用いた音声認識方法の詳細は、中川聖一著「確率モデルによる音声認識」(電子情報通信学会)に詳しい。このHMMに基づく話者照合や話者適応や話者正規化に関する研究が行われている。通常、話者正規化や話者適応技術は、音声データの内容や量に依存するため、少量の発声データからでは安定した性能向上が難しい。そこで、声道長を用いた手法が注目されており、特に声道長に基づく話者正規化が盛んに研究されて効果が出ている。

【0003】上記声道長は、音声のスペクトルの大まかな特徴を表わすパラメータである。そして、上記声道長の差は話者間の主な変動要因であり、声道長は従来の話者適応法に比べて1個のパラメータあるいは極めて少ないパラメータで音声の特徴を表現できることから、声道長にはより少量の学習データで効率良く正規化できるというメリットがある。

【0004】ところで、標準話者の音声パターンに対する入力話者の音声サンプルの尤度を最大にするという基準(最尤推定)に従って、上記音声サンプルにおける周波数軸の線形伸縮係数 α (声道長正規化係数)を求める

(ML-VTLN法: Maximum Likelihood Vocal Tract Length Normalization)。そして、この声道長伸縮係数 α を用いて、入力話者の音声サンプルの周波数軸を線形伸縮して話者正規化を行う技術が提案されている(例えば、AT&T Bell Labs. Li Lee, Richard C. Rose, "Speaker Normalization using Efficient Frequency Warping Procedures", pp. 353-356 ICASSP96 (1996))。また、特開平11-327592号公報においては、声道を前室と後室との2つの室に分け、入力音声のフォルマント周波数を用いて、各室に対応した2つの周波数軸線形伸縮係数 α を用いて話者正規化する技術が開示されている。

【0005】尚、上記話者適応は標準となる音響モデルを入力話者に対して適応(つまり正規化)させる技術であり、話者正規化とは表裏一体の関係にある。

【0006】また、話者クラスタリングを用いた音声認識方法がある。この音声認識方法においては、学習話者間の距離を定義して学習話者をクラスタリングしておき、クラス毎にそのクラスに属する学習話者群の音声データを用いて音響モデルを作成する。そして、認識時には、入力音声に最適なクラスを選択し、そのクラスの音響モデルを用いて認識処理を行うのである。その場合における学習話者間の距離として上記声道長の周波数軸線形伸縮係数を用いる音声認識装置が提案されている(特開平11-175090号公報)。この公報にお

いては、声道を前室と後室との2つの室に分け、各室に対応した2つの周波数軸線形伸縮係数を用いて学習話者をクラスタリングするようにしている。

【0007】また、声道長の非線形な伸縮関数を導入してその係数 α でクラスタリングする方法や、GMM(ガウシアン混合モデル)を用いて話者クラスタリングする方法が提案されている(佐藤他「GMMによる音響モデル用学習データの自動選択」日本音響学会春季研究発表会講演番号2-8-3 2000年3月)。上記GMMは1状態の混合ガウス分布で表現される音響モデルであり、発話内容に因らずに入力音声に声質の近いGMMが大きい値を出力するように設計されている。元々は話者照合における話者モデルとして提案された手法である。

【0008】

【発明が解決しようとする課題】しかしながら、上記従来の声道長に基づく話者適応や話者正規化には、以下のような問題がある。すなわち、声道長伸縮関数の求め方として学習サンプル全体を対象として最尤推定する方法(ML-VTLN法)等が提案されている。このような声道長に基づく話者適応や話者正規化は極めて少ないパラメータ数で表現できるとは言うものの、声道長の抽出は発声データの内容や量に大きく左右されるために、少ない学習サンプルから必ずしも安定して声道長を抽出できるとは限らない。したがって、声道長に基づいて話者正規化や話者適応や話者クラスタリングを行う音声認識装置においては、性能劣化を招くと言う問題がある。

【0009】実際の声道長はMRI(磁気共鳴画像診断装置)で測定しなければ分からないため、現時点においては直ちに真の声道長を知るのは困難な状況にある。上記特開平11-327592号公報および特開平11-175090号公報では、声道パラメータを得るために入力音声のフォルマント周波数を用いている。しかしながら、一般的にフォルマント周波数を全自動で求めることは困難であり、上記特開平11-327592号公報に開示された線形伸縮係数を用いた話者正規化方法や上記特開平11-175090号公報に開示された線形伸縮係数を用いた音声認識装置では、実時間性に欠けると言う問題がある。

【0010】また、上記特開平11-175090号公報のごとく、話者クラスタリングを用いた音声認識のアプローチも盛んに試みられているが、大きな性能改善は達成できていない。不特定話者(SI)音響モデル(すなわち男女共通の音響モデル)をベースラインとすると、男女別(GD)音響モデルは最もシンプルながら性能向上量が最も大きい。しかしながら、話者クラスによって更なる細分化(クラスタ化)を行っても効果は薄いという報告がなされており、その場合における単語誤り率(WER: Word Error Rate)の削減は10%~20%程度に留まっている。これは、話者間の距離を定義する適当な尺度がないために上手くクラスタリングできなかった

り、クラスタを増やすと1つのクラスタ当りの学習話者数が少なくなってロバスト性に欠けたりするためである。

【0011】さらに、何れの音響モデルの場合も、各話者クラスタの境界領域では学習サンプルが希薄だったり段差がでたりしているため上手く学習されていない。したがって、入力話者が各クラスタの境界付近に位置する場合には、認識率が劣化するという問題(所謂、hard decision問題)が生ずることになる。尚、個々の学習話者の音響モデル間の距離でクラスタリングを行った場合は、クラスタを木構造にし、入力話者が二つのクラスタの境界付近に位置する場合は上記2つのクラスタの上位ノードのクラスタの音響モデルを採用する方法もある。しかしながら、この方法の場合には、二つのクラスタの境界付近に位置する入力話者に対しては上位ノードの音響モデルを使用するためによりブロードな音響モデルとなってしまう、高い認識率は得にくいのである。

【0012】ところで、上記ML-VTLN法に基づいて話者をクラスタリングする場合には、以下のような問題がある。

- ・真の声道長伸縮係数 α の値を求めるのは困難である。上記真の声道長伸縮係数 α 値を求めるには各話者についてMRI装置で実測しなければならない。しかしながら、既に構築済みの学習用音声データベースがあり、直ちにはそのデータベースを活用するしかない場合や、上記MRI装置を利用し難い環境下にある場合には、音声波形から声道長伸縮係数 α の値を自動推定する必要がある。したがって、自動推定する限りにおいてはどうしても推定誤差の問題が付きまとうことになる。

- ・例え、上記MRI装置で測定した実測値をもってしても、発声の仕方の影響があるために、適切な声道長伸縮係数 α の値が得られるとは限らない。

【0013】一方、上記GMMに基づいて話者をクラスタリングする場合には、一般に以下のような課題がある。

- ・初期値をランダムにして全自動でクラスタリングし、その後はHMMの学習アルゴリズムに頼っている。しかしながら、音声データは多数の要因が絡み合って複雑な構造を成しているために、このような方法の場合には、音声の微細な特徴を捉えてクラスタリングする危険性がある。

- ・上記GMM間の距離の物理的意味が不明である。つまり、距離の大小が音響的に何に対応しているのかが分からないために、周波数伸縮による話者正規化は適用できない。

- ・クラスタ化による学習データの減少を補う目的で近傍クラスタの学習データを編入させる場合に上記GMM間の距離を用いると、話者の特徴空間上、様々な方向に位置する話者データを編入するになる。その結果、ばやけた分布になってしまう、精密な話者特徴を抽出できなく

なる。したがって、このようにしてできたGMMを基に学習されたHMMに対しても精度の劣化を招くと言う問題がある。

【0014】以上のごとく、上記話者適応(話者正規化)においては少ない発声データから音響モデルを精度良く適応できないため、誤り率を半減させるためには数十単語以上の発声データが必要となり、学習話者に負担を強いることになるという問題がある。また、音声合成における声質変換の場合にも、同様に少ない発声データからは精度良く声質が得られないという問題がある。

10 【0015】そこで、この発明の目的は、より少ない発声データから精度良く話者特徴を抽出できる話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体を提供することにある。

【0016】

【課題を解決するための手段】上記目的を達成するため、第1の発明は、入力話者の音声に基づいて、標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として抽出する話者特徴抽出装置において、各学習話者に関して、上記標準話者に対する声道長の伸縮係数 α を所定の方法によって予め求め、この求められた伸縮係数 α の値に基づいて上記学習話者をクラスタリングする学習話者クラスタリング手段と、上記クラスタリングされた各クラスタに属する話者集合毎に、学習によってGMMを生成する音響モデル生成手段と、上記生成されたGMMの夫々に対する上記学習話者の音声サンプルの尤度を算出し、その尤度に基づいて上記学習話者を再クラスタリングする再クラスタリング手段と、上記音響モデル生成手段と再クラスタリング手段とを制御して、所定の条件を満たすまで、上記GMMの生成と上記学習話者の再クラスタリングとを繰り返すループ学習手段と、上記ループ学習手段によって最終的に生成されたGMMの群を格納する音響モデル格納部と、上記音響モデル格納部に格納されたGMMの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈するGMMを入力話者が属するクラスタの情報として選択する話者クラスタ選択部を備えて、上記入力話者の特徴として上記最大の尤度を呈するGMMを抽出することを特徴としている。

【0017】上記構成によれば、学習話者をクラスタリングするに際して、先ず、各学習話者を標準話者に対する声道長の伸縮係数 α に基づいてクラスタリングし、各クラスタに属する話者集合毎にGMMを生成する。そして、このGMM群を用いてループ学習を行うことによって学習話者を再クラスタリングするようにしている。こうして、各クラスタの初期値として声道長という大局的な特徴を明示的に与えることによって、各クラスタ間の距離の物理的意味が明確になり、効率よくクラスタリングが行われる。さらに、1状態の混合ガウス分布型音響

モデルであるGMMを用いて学習話者を再クラスタリングすることによって、発話内容に因らずに話者の特徴を良く表わす話者クラスタが得られ、声道長伸縮係数 α の抽出誤りも修復されている。

【0018】したがって、上述のようにして得られた話者クラスタ毎にGMMが格納された音響モデル格納部を用いて、話者クラスタ選択部によって、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択することによって、発話内容に因らずに精度良く入力話者の特徴が抽出される。

【0019】また、1実施例では、上記第1の発明の話者特徴抽出装置において、上記学習話者クラスタリング手段によってクラスタリングされた各クラスタのうちの注目クラスタに隣接する隣接クラスタに属する学習話者の音声サンプル、または、上記注目クラスタと上記伸縮係数 α 値の差が所定値以内の近傍クラスタに属する学習話者の音声サンプルに対して、上記注目クラスタと隣接クラスタまたは近傍クラスタとの上記伸縮係数 α 値に基づいて周波数伸縮を行うことによって上記注目クラスタに属する音声サンプルを生成し、この生成された音声サン

プルを上記注目クラスタに編入して当該クラスタの音声サンプル数を豊富化する操作を、上記学習話者クラスタリング手段によってクラスタリングされた総てのクラスタについて実行する音声サンプル豊富化手段を備えて、上記音響モデル生成手段は、上記音声サンプル豊富化手段によって音声サンプル数が豊富化された後の各クラスタ毎に、上記GMMを生成するようにしている。

【0020】この実施例によれば、上記学習話者のクラスタリングに際して、ある注目クラスタの隣接クラスタまたは近傍クラスタに属する話者の音声サンプルに対して周波数伸縮が行われ、上記注目クラスタに属する音声サンプルが生成されて注目クラスタに編入される。こうして、学習データ不足が補われて、少ない発声データからでも各クラスタの音響モデルが精密に構築される。

【0021】また、1実施例では、上記第1の発明の話者特徴抽出装置において、上記注目クラスタに属する音声サンプルを生成する際に、上記音声サンプル豊富化手段が上記学習話者の音声サンプルに対して周波数伸縮を行う音声区間を、有音・無音の別および調音点に基づいて限定するようにしている。

【0022】この実施例によれば、上記音声サンプル豊富化手段によって、上記隣接クラスタや近傍クラスタの音声サンプルから注目クラスタに属する音声サンプルを生成する際に、上記音声サンプルに対して周波数伸縮を行う音声区間が有音・無音の別および調音点に基づいて限定される。したがって、声道長の差の影響を受け難い音素や無音部を上記周波数伸縮の対象外にして、声道長の差の影響を受け難い音素や無音部まで変形されることが防止される。

【0023】また、1実施例では、上記第1の発明の話

者特徴抽出装置において、上記再クラスタリング手段によって上記学習話者を再クラスタリングする場合に、上記再クラスタリングの対象となる対象学習話者が再クラスタリングの前に属していたクラスタの伸縮係数 α と再クラスタリング後に属するクラスタの伸縮係数 α とが所定値以上離れている場合には、当該対象学習話者を上記再クラスタリングの対象から外すようになっている。

【0024】この実施例によれば、再クラスタリングの対象となる学習話者が再クラスタリングの前後に属しているクラスタの伸縮係数 α が所定値以上離れている場合は、当該対象学習話者を上記再クラスタリングの対象から外すことによって、声道長伸縮係数 α が極端に異なる話者同士が同じクラスタに属することが防止される。

【0025】また、1実施例では、上記第1の発明の話者特徴抽出装置において、上記ループ学習手段によって最終的にクラスタリングされた各クラスタに属する学習話者を更にクラスタリングしてサブクラスタを生成し、上記各サブクラスタに属する話者集合毎に学習によってGMMを生成するサブクラスタ生成手段を備えて、上記音響モデル格納部は、上記サブクラスタ生成手段によって生成されたGMMの群を、各サブクラスタの伸縮係数 α に対応付けて格納するようになっている。

【0026】この実施例によれば、上記ループ学習手段によって最終的にクラスタリングされた各クラスタがさらにクラスタリングされてサブクラスタが生成される。このサブクラスタは声道長以外の要因にも対応することができ、より精密な話者特徴が抽出される。

【0027】また、第2の発明は、入力話者の音声に基づいて標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として抽出する話者特徴抽出装置において、上記標準話者に対する声道長の伸縮係数 α の値に基づいて学習話者をクラスタリングし、各クラスタに属する話者集合毎にGMMを生成し、この生成されたGMM夫々に対する上記学習話者の音声サンプルの尤度に基づいて上記学習話者を再クラスタリングし、上記GMMの生成と上記学習話者の再クラスタリングとを所定の条件を満たすまで繰り返して最終的に生成されたGMMの群が格納された音響モデル格納部と、上記音響モデル格納部に格納されたGMMの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈するGMMを入力話者が属するクラスタの情報として選択する話者クラスタ選択部を備えて、上記入力話者の特徴として上記最大の尤度を呈するGMMを抽出することを特徴としている。

【0028】上記構成によれば、標準話者に対する声道長の伸縮係数 α に基づいて学習話者をクラスタリングし、各クラスタ毎のGMMの生成とそのGMM群を用いた学習話者の再クラスタリングとを所定の条件を満たすまで繰り返し、最終的に生成されたGMM群を話者クラスタ毎に格納した音響モデル格納部を用いて、話者クラ

スタ選択部によって、入力話者の音声サンプルに対して最大尤度を呈するGMMが選択される。こうして、発話内容に拘らず精度良く入力話者の特徴が抽出される。

【0029】また、第3の発明は、音響モデルとしてHMMを用い、入力話者の音声に基づいて抽出された標準話者の音声と上記入力話者の音声との関係を表わすパラメータを話者特徴として上記入力話者の音声を認識する音声認識装置であって、上記第1の発明または第2の発明の話者特徴抽出装置と、上記話者特徴抽出装置の音響モデル格納部に格納された各GMMによって表わされる話者クラスタに属する話者集合毎に、学習によって生成されたHMMの群を格納するHMM格納部と、上記話者特徴抽出装置によって選択されたクラスタに基づいて、上記HMM格納部に格納されている上記選択されたクラスタに対応するクラスタのHMMを音声認識用の音響モデルとして切り換え選出する切換部を備えたことを特徴としている。

【0030】上記構成によれば、上記第1の発明または第2の発明の話者特徴抽出装置の音響モデル格納部における話者クラスタ毎に生成されたHMMの群が格納されたHMM格納部から、切換部によって、上記話者特徴抽出装置で選択された話者クラスタのHMMが音声認識用の音響モデルとして切り換え選出される。こうして、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わすHMMを用いて、入力話者の音声

が正確に認識される。

【0031】また、第4の発明は、音響モデルとしてHMMを用い、入力話者の音声に基づいて抽出された標準話者の音声と上記入力話者の音声との関係を表すパラメータを話者特徴として上記入力話者の音声を認識する音声認識装置であって、上記第1の発明あるいは第2の発明に係る音響モデル格納部に格納された各GMMによって表わされる話者クラスタに属する話者集合毎に、学習によって生成されたHMMの群を格納するHMM格納部と、上記HMM格納部に格納されたHMMの夫々に対する入力話者の音声サンプルの尤度を算出し最大の尤度を呈するHMMを音声認識用の音響モデルとして切り換え選出する切換部を備えたことを特徴としている。

【0032】上記構成によれば、上記第1の発明または第2の発明に係る音響モデル格納部の話者クラスタ毎に生成されたHMMの群が格納されたHMM格納部を用いて、切換部によって、最大の尤度を呈するHMMが音声認識用の音響モデルとして切り換え選出される。こうして、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わすHMMを用いて、入力話者の音声

が正確に認識される。

【0033】また、第5の発明は、声道長の伸縮関数 α を用いて入力音声のスペクトルの周波数軸を伸縮することによって入力話者の音響特徴量を標準話者の音響特徴量に正規化する話者正規化手段を有する音声認識装置に

おいて、上記話者正規化手段は、上記第1の発明あるいは第2の発明の話者特徴抽出装置と、上記入力話者の音声サンプルに基づいて上記話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α を用いて、上記入力音声のスペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【0034】上記構成によれば、入力話者の音声サンプルに基づいて上記第1の発明あるいは第2の発明の話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α を用いて、周波数ワープ手段によって上記入力音声のスペクトルの周波数軸が伸縮されて、上記入力話者の音響特徴量が標準話者の音響特徴量に正規化される。こうして、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わす声道長の伸縮係数 α を用いて、より標準話者の音響特徴量に近づくように話者正規化が行われる。その結果、高い音声認識率が得られる。

【0035】また、第6の発明は、声道長の伸縮関数 α を用いて音声のスペクトルの周波数軸を伸縮することによって音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、上記話者適応手段は、上記第1の発明あるいは第2の発明の話者特徴抽出装置と、上記入力話者の音声サンプルに基づいて上記話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α の逆数を用いて、上記音響モデルのスペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【0036】上記構成によれば、入力話者の音声サンプルに基づいて、上記第1の発明あるいは第2の発明の話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α の逆数を用いて、周波数ワープ手段によって音響モデルのスペクトルの周波数軸が伸縮されて上記音響モデルが入力話者に話者適応される。こうして、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わす声道長の伸縮係数 α の逆数を用いて、より入力話者の音響特徴量に近づくように話者適応が行われる。その結果、高い音声認識率が得られる。

【0037】また、第7の発明は、入力話者の音声に基づいて、標準話者の音声と上記入力話者の音声との関係を表すパラメータを話者特徴として抽出する話者特徴抽出方法であって、各学習話者に関して、上記標準話者に対する声道長の伸縮係数 α を所定の方法によって予め求め、この求められた伸縮係数 α の値に基づいて上記学習話者をクラスタリングし、上記クラスタリングされた各クラスタに属する話者集合毎に学習によってGMMを生成し、上記生成されたGMMの夫々に対する上記学習話者の音声サンプルの尤度を算出し、その尤度に基づいて上記学習話者を再クラスタリングし、所定の条件を満た

すまで上記GMMの生成と上記学習話者の再クラスタリングとを繰り返すループ学習を行い、上記ループ学習によって最終的に生成されたGMMの群を音響モデル格納部に格納し、上記音響モデル格納部に格納されたGMMの夫々に対する入力話者の音声サンプルの尤度を算出し、最大の尤度を呈するGMMを上記入力話者の特徴として抽出することを特徴としている。

【0038】上記構成によれば、上記第1の発明の場合と同様に、学習話者をクラスタリングする際に、各クラスタの初期値として声道長という大局的な特徴を明示的に与えることによって、各クラスタ間の距離の物理的意味が明確になり、効率よくクラスタリングが行われる。さらに、1状態の混合ガウス分布型音響モデルであるGMMを用いて学習話者を再クラスタリングすることによって、発話内容に因らずに話者の特徴を良く表わすクラスタが得られ、声道長伸縮係数 α の抽出誤りも修復されている。

【0039】したがって、上述のようにして得られた話者クラスタ毎にGMMが格納された音響モデル格納部を用いて、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択することによって、発話内容に因らずに精度良く入力話者の特徴が抽出される。

【0040】また、第8の発明のプログラム記録媒体は、コンピュータを、上記第1の発明に係る学習話者クラスタリング手段、音響モデル生成手段、再クラスタリング手段、ループ学習手段、音響モデル格納部および話者クラスタ選択部として機能させる話者特徴抽出処理プログラムが記録されていることを特徴としている。

【0041】上記構成によれば、上記第1の発明の場合と同様に、発話内容に因らずに話者の特徴を良く表わす話者クラスタ毎にGMMが格納された音響モデル格納部を用いて、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択することによって、発話内容に因らずに精度良く入力話者の特徴が抽出される。

【0042】

【発明の実施の形態】以下、この発明を図示の実施の形態により詳細に説明する。

<第1実施の形態>図1は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者クラスタリング方式を用いた音声認識装置である。音声入力部1において、マイクから入力された音声はデジタル波形に変換されて音響分析部2に入力される。音響分析部2は、入力されたデジタル波形を短い時間間隔(フレーム)毎に周波数分析し、スペクトルを表す音響パラメータのベクトル系列に変換する。ここで、上記周波数分析には、MFCC(メル周波数FFT(高速フーリエ変換)ケプストラム)やLPC(線形予測分析)メルケプストラム等のスペクトルを効率よく表現できる方法が用いられる。こうして得られた音響パラメータ系列は、話者クラスタ選択部3及び尤度(音韻類似度)演算部

4に送出される。

【0043】上記話者クラスタ選択部3は、GMM格納部6と共に話者特徴抽出部1を構成し、以下のようにして話者特徴としてのクラスタ情報を生成する。すなわち、話者クラスタ選択部3は、入力された音響パラメータ系列にGMM格納部6に話者クラスタ別に格納された n 個のGMMの夫々を作用させて尤度を算出する。そして、算出された n 個の尤度のうちの最大値を与えるGMMのインデックス(i)($i=1, 2, \dots, n$)を、その入力話者に適合した話者クラスタ情報として出力する。ここで、上記GMMは、1状態からなる混合ガウス分布で表現される。

【0044】切換部5は、話者クラスタ音響モデル格納部7に話者クラスタ別に格納された音響モデル(本実施の形態ではHMMを使用)の中から、話者クラスタ選択部3からのクラスタ情報に適合する話者クラスタの音響モデルを切り換え選択して尤度演算部4に送出する。そうすると、尤度演算部4は、音響分析部2からの入力音声の音響パラメータベクトルに対して切換部5からの音響モデルを作用させて、各音韻の状態毎に尤度を算出する。そして、得られた尤度系列を照合部8に送出する。

【0045】上記照合部8は、上記尤度演算部4からの尤度系列に対して、言語モデル格納部9に登録された総ての言語モデル(単語)との照合を行ない、各単語のスコアを算出する。そして、上位のスコアを呈する単語を認識候補(認識結果)として出力部10から出力するのである。

【0046】ここで、上記話者クラスタ選択の方法には、以下の[a]および[b]に示す2通りの方法がある。本実施の形態においては[a]の方法を用いている。

[a] GMM格納部6の利用

[b] 話者クラスタ音響モデル格納部7の利用

【0047】上記[a]の方法は、各話者クラスタに対して1つのGMMを作成しておく。そして、入力音声に対して各GMMを作用させてGMM毎の尤度を算出し、最も大きい尤度を算出したGMMに対応する話者クラスタを選択するのである。その場合、入力音声の正解音素列をユーザが教える必要がなく、教師なしでクラスタが選択できる。すなわち、エンロールモードがないシステムにおいて有効なのである。

【0048】また、上記[b]の方法は、上記話者クラスタの音響モデル自身の尤度を用いる方法である。入力音声に対して教師語彙が与えられ、各話者クラスタにおける教師語彙の音響モデルを用いて認識処理を行って、話者クラスタ毎の尤度を算出する。そして、最も大きい尤度を呈する話者クラスタを選択するのである。この選択方法は、エンロールモードにおいて入力音声の正解音素列をユーザが教えるという教師あり学習を基本としている。認識処理と同じ高精度な音響モデルを用いるので計算量は多くなるが、エンロールによって正確なクラスタ

選択が可能となる。

【0049】ここで、本実施の形態における話者クラス
タ音響モデル格納部7は、学習話者のクラス数 n に応
じて、第1音響モデル格納部7a, 第2音響モデル格納部
7b, ..., 第 n 音響モデル格納部7nの n 個の音響モデル格
納部で構成されている。ここで、各音響モデル格納部7
a~7nに格納される各音響モデルは、混合ガウス分布型
のHMMである。この発明においては、生理的な特徴の
変動に対処可能にすることを目的としており、話者性の
大局的な安定要因である声道長を初期値としたGMMを
クラスタリング対象にするのである。

【0050】尚、生理的な特徴の変動要因としては、上
記声道長以外にも鼻腔、副鼻腔、声帯等の多くの要因が
あり、それらが絡み合って複雑な特徴を成している。した
がって、個々の要因を数理的に扱うのは得策ではない。
そこで、本実施の形態においては、各要因の複雑な特徴
を混合ガウス分布型HMM(音響モデル)や、同一クラス
タ内での複数のGMMによるサブクラスタで表現するの
である。

【0051】以下、上記GMM格納部6に格納されるG
MM群の作成方法について詳細に説明する。本実施の形
態におけるGMM群の作成方法は、下記の処理手順によ
って行われる。これらの処理はGMM学習部12によ
って、学習データ格納部13に格納された学習話者の音声
データを用いて行われる。尚、記憶領域や処理量が膨大
になるために、予めオフライン処理によって作成されて
いる。作成されたGMM群は、通常はROM(リード・オ
ンリ・メモリ)やフラッシュメモリやハードディスク等で
構成されるGMM格納部6に格納される。尚、学習デー
タ格納部13に格納された全学習話者の集合を男女別に
2分割してもよい。その場合には、男性用話者クラスタ
と女性用話者クラスタとの合計2種類のGMM群が生成
されることになる。

【0052】① 夫々の学習話者の音声データに所定の
方法によって声道長伸縮係数 α を与え、声道長伸縮係数
 α の値に応じて学習話者を N 個にクラスタリングしてお
く。尚、GMMの初期モデルの混合数は M とする。

② あるクラスタCに隣接するクラスタに属する音声デ
ータ、または、あるクラスタCとの声道長伸縮係数 α の
値の差が所定値 δ 以内であるクラスタに属する音声デー
タに対して、ある特定の区間を対象とした周波数伸縮を
行うことによって上記クラスタCに属する音声データを
生成する。そして、この生成された音声データをクラス
タCに編入する。こうして、クラスタCの音声データを
豊富化するのである。この豊富化処理を総てのクラスタ
について行う。

③ GMMパラメータを各クラスタ内の音声データから
ML(最尤)アルゴリズムにより推定することによって、
各クラスタのGMMを作成する。

④ 作成されたGMMに対する各学習話者の音声データ

のフレーム平均尤度を算出する。

⑤ ある話者の音声データに対して最も高いフレーム平
均尤度を与えるGMMのクラスタに、その話者を移動さ
せる(編入する)。但し、そのクラスタの声道長伸縮係数
 α 値に比べて、声道長伸縮係数 α 値が所定値 ε 以上離れ
ている話者については、そのクラスタには編入しない。

⑥ 移動させる話者がなくなるか、予め設定した最大の
繰り返し回数になるまで③~⑤の処理を繰り返す。

⑦ 混合数を1つ増加して③~⑥の処理を行う。

⑧ 所望の混合数になるまで③~⑦の処理を繰り返す。

⑨ ①~⑧で得られた学習話者のクラスタ結果を基に、
サブクラスタ化したGMMを作成する(オプション)。

【0053】尚、上記GMM群の作成処理手順①におけ
る各学習話者に対する声道長伸縮係数 α の付与は、ML
-VTLN法や、広母音の第2フォルマントの存在領域
以下の領域を部分的に補正した非線形周波数ワーピング
関数を用いたML法等によって与えられる。尚、声道長
伸縮係数 α の推定に際しては、後述する音声区間の分類
を用いてもよい。また、MRI装置を利用できる等、各
学習話者の声道長を実測できる環境にある場合には、実
測された α を用いてもよい。さらに、クラスタ数 N と混
合数の初期値 M とは、学習データ量や声道長伸縮係数 α
の信頼性に依存するが、例えば $N=12$, $M=20$ 等と
する。 $N=12$ とは、 α 軸上の区間(0.88, 1.12)
を0.02刻みに分割することに相当する。また、各ク
ラスタの境界をオーバーラップするような分割を許して
もよい。

【0054】また、上記GMM群の作成処理手順②にお
ける音声データの変換の際には、線形周波数伸縮関数
や、広母音の第2フォルマントの存在領域以下の領域を
部分的に補正した折れ線周波数伸縮関数を用いる。以
下、②の処理内容について詳しく述べる。あるクラスタ
Cに隣接するクラスタ、または、あるクラスタCとの声
道長伸縮係数 α の値の差が所定値 δ 以内であるクラスタ
に属する話者の音声データに対して、声道長伸縮係数 α
値に基づいて周波数伸縮を行うことによって、クラスタ
Cに属する音声データを生成するのである。例えば、 α
 $=1.05$ のクラスタDに属する音声データに基づいて
 $\alpha=1.03$ のクラスタCに属する音声データを生成す
る際には、周波数を0.98だけ伸縮する。そして、生
成された学習話者の音声データをクラスタCに編入する
のである。尚、生成前の音声データは元のクラスタDに
属したままにしておく。但し、学習話者数および1話者
当りのデータ量が大量にある場合は、この処理は省略し
ても差し支えない。逆に、学習話者数や1話者当りのデ
ータ量が少ない場合には、上記所定値 δ を大きく(例え
ば $\delta=0.05$)に設定して編入させるデータを増やす。

【0055】上記GMM群の作成処理手順②における音
声データの変換の際に、周波数伸縮を行う対象となる音
声区間の分類については後述する。尚、伸縮量が非常に

小さい場合には全区間を周波数伸縮対象としてもよい。

【0056】上記GMM群の作成処理手順⑤において、上記所定値 ε の値は、①で与えられた声道長伸縮係数 α の信頼度に依存する。上記MRI装置によって実測した場合のように声道長伸縮係数 α の信頼度が高い場合は移動を禁止してもよいし、所定値 ε を小さい値に設定してもよい(例えば $\varepsilon = 0.02$)。逆に、声道長伸縮係数 α の信頼度が低い場合には、所定値 ε の値を大きくする(例えば $\varepsilon = 0.04$)。ところで、上記移動の際に、通常は周波数伸縮(α の値を書き換えることに相当)を行わ

ない。但し、所定値 ε の値を大きめに設定した場合には周波数伸縮を行ってもよい。

【0057】上記GMM群の作成処理手順⑨におけるサブクラス化はオプションであり、処理方法は後述する。

【0058】次に、上記GMM群の作成処理手順①にお

分類	音素	説明
[a]	無音区間	環境騒音そのものであり、声道長とは全く関係ない。
[b]	調音点が歯茎より前に位置する子音	音源が声道出口付近に位置するために、声道全体による共鳴管が形成されず、声道長の影響を受け難い。
[c]	調音点が歯茎より後に位置する子音、半母音	母音と同様に声道長に直接影響される。しかしながら、子音や半母音区間は元来安定していないので声道長の推定には適さない。
[d]	「ウ」を除く母音	母音は全般的に声道長の影響を直接受ける。
[e]	母音「ウ」、撥音	日本語の「ウ」は発声の仕方によってフォルマンと周波数が大きく変動するので、声道長の影響よりも発声の仕方に大きく影響される。撥音も音素環境に大きく依存すると共に鼻音化の影響が大きい。

【0060】そして、この分類に基づいて、以下のよう

な区別に従って、推定/正規化時におけるGMM学習部

12の処理を制御するのである。
・上記GMM群の作成処理手順①における声道長伸縮係数 α の初期値推定時…分類[d]
・上記GMM群の作成処理手順②における正規化時…分類[c], 分類[d], 分類[e], (分類[b])
但し、上記正規化時には分類[b]を含めてもよい。発音の仕方によっては、音素「イ」も音素「ウ」と同様に挟母音なのでフォルマンと周波数が大きく変動する場合がある。したがって、分類[e]に音素「イ」を含め、分類[d]から音素「イ」を除いてもよい。

【0061】次に、上記GMM群の作成処理手順⑨にお

けるサブクラス化の作成方法について説明する。サブクラスの作成は、上記GMM群の作成処理手順①～⑧によって得られた総てのクラスに対して、下記の処理手順を繰り返すことによって行う。

【0062】A) クラス内の学習話者をランダムにL

個に分割する。すなわち、L個のサブクラスを与える

のである。但し、GMM初期モデルの混合数は1とす

る。

B) GMMパラメータを各サブクラス内の音声データ

からMLアルゴリズムにより推定することによって、各

*ける声道長伸縮係数 α の推定や、②における周波数伸縮の際に、対象とする音声区間の分類について説明する。

上述の例において、クラスDに属している音声データをクラスCに変換するということはクラスCに正規化していることを意味しているので、ここでは、②の周波数伸縮を単に正規化と呼ぶことにする。先ず、入力話者の音声データに、不特定話者用音響モデルまたは選択された話者クラス音響モデルまたは特定話者音響モデルから選択されたものを用いたビタビアルゴリズムによって、音素境界情報を求めておく。

【0059】次に、その境界情報に基づいて、上記音声データのうち推定/正規化処理の対象となる区間を制御するのである。本実施の形態においては、この推定/正規化処理の対象となる区間を、表1に示す5種類に分類する。

表1

サブクラスのGMMを作成する。

C) 作成されたGMMに対する各学習話者の音声データのフレーム平均尤度を算出する。

D) ある話者の音声データに対して最も高いフレーム平均尤度を与えるGMMのサブクラスに、その話者を移動させる。

E) 移動させる話者がなくなるか、予め設定した最大の繰り返し回数になるまでB)～D)の処理を繰り返す。

F) 混合数を1つ増加してB)～E)の処理を行う。

G) 所望の混合数になるまでB)～F)の処理を繰り返す。

【0063】上記サブクラスの作成処理手順A)におけるサブクラス数Lの値は、クラス内の学習話者数及び1話者当りのデータ量に依存するが、通常2～10の間に設定する。クラスによってサブクラス数Lの値を変えてもよい。例えば、分布の中心である $\alpha = 1.0$ 付近のクラスは話者数が多いため $L = 5$ とする。一方、分布の周辺である $\alpha = 0.9$ 付近や $\alpha = 1.1$ 付近のクラスは話者数が少ないため $L = 2$ 等にするのである。また、上記サブクラスの作成処理手順G)における所望の混合数も、サブクラス内の学習話者数および1話者当りのデータ量に依存させてもよい。

【0064】尚、本実施の形態における話者クラスタリ

ングの場合や第2実施の形態における話者正規化の場合での α は、入力音声から標準音声への正規化係数である。これに対して、第3実施の形態における話者適応の場合での α は、標準音声から入力音声への写像係数である。このように、両者は裏表の関係であるため、 α の値は逆数の関係になる。

【0065】次に、上記話者クラスタ音響モデル格納部7に格納されるHMM群の作成方法について説明する。上記GMM作成時においてクラスタリングされた話者クラスタ毎にMLアルゴリズムを用いて学習することによって、混合ガウス分布型のHMMが作成される。尚、上記GMM群の作成処理手順⑨においてサブクラスタ化が行われている場合には、サブクラスタ毎に、同様の手法によって混合ガウス分布型のHMMが作成される。これらの処理はHMM学習部14で行われるのであるが、記憶領域や処理量が膨大になるため予めオフライン処理によって作成されている。そして、作成されたHMM群は、通常はROMやフラッシュメモリやハードディスク等で構成される話者クラスタ音響モデル格納部7に格納される。

【0066】上記構成において、入力音声の認識時には以下のように動作する。先ず、話者クラスタ選択部3によって、上述のようにして最適な話者クラスタが選択され、クラスタ情報(i)が切換部5に送出される。次に、尤度演算部4によって、上記切換部5で切り換え選択された話者クラスタの音響モデルを用いて尤度演算が行われ、得られた尤度系列が照合部8に送出される。そして、照合部8によって、ビタビサーチ等の探索アルゴリズムが用いられて言語モデル格納部9の言語モデルとの照合が行われ、各単語のスコアが算出される。尚、本実施の形態においては、照合部8による照合処理の前段処理が訴求点であるから、照合部8に関する詳細な説明は省略する。

【0067】上述したように、上記GMMは1状態の混合ガウス分布で表現される音響モデルであって、発話内容に因らずに入力音声に声質の近いGMMが大きい値を出力するように設計されており、話者の特徴を表わすには好適である。ところが、GMM間の距離の物理的意味が不明であるため、GMM間の距離でクラスタ化を行った場合にはばやけた分布になってしまい、精密な話者特徴を抽出できないという問題がある。

【0068】そこで、本実施の形態においては、GMM学習部12で学習話者をクラスタリングするに際して、先ず、学習データ格納部13に格納された各学習話者の音声データに、声道長正規化手法によって求められた声道長伸縮(正規化)係数 α の値を与える。そして、この声道長伸縮係数 α に基づいて学習話者をクラスタリングし、各クラスタに属する話者集合毎に所定の学習方法によって1状態の混合ガウス分布型音響モデルであるGMMを生成する。そして、次に、このGMM群を用いてル

ープ学習を行って、学習話者をクラスタリングし直すのである。

【0069】音声データは多数な要因がからみあって複雑な構造を成しているため、GMM間の距離でクラスタ化を行った場合には音声の微細な特徴を捉えてクラスタリングしてしまう危険性がある。そこで、上述のように、声道長という大局的な特徴を明示的に与えることによって、距離の物理的意味が明確になって、以後の学習をスムーズに実行でき、より効率よくクラスタリングできるのである。さらに、各クラスタのGMMを生成し、このGMM群を用いて学習話者をクラスタリングし直すようにしている。したがって、発話内容に因らずに話者の特徴を良く表わすクラスタを得ることができ、声道長伸縮係数 α の抽出誤りも修復できるのである。

【0070】また、上記学習話者のクラスタリングに際して、あるクラスタCとの声道長伸縮係数 α の差が δ 以内である近傍のクラスタDに属する話者の音声データに基づいて、声道長正規化手法によって上記クラスタCに属する音声データを生成するようにしている。したがって、話者をクラスタリングする際における学習データ不足を補うことができ、各話者クラスタの音響モデルを精密に構築できるのである。このことは、音声認識時におけるクラスタの選択(声道長の推定)をより正確に且つ安定して行うことができ、認識率の向上につながる。また、分布の周辺に位置する話者の認識率を向上させることができる。

【0071】また、上述のようにして上記声道長に基づいて求められた1つの話者クラスタを複数のサブクラスタに分割し、夫々のサブクラスタのGMMを生成してこれを話者特徴とする。こうして、生成されたサブクラスタは声道長以外の要因にも対応することができるので、上記サブクラスタを用いることによってより精密な話者特徴を抽出することができるのである。

【0072】したがって、上述のようにしてクラスタリングされた話者クラスタ毎に作成されたHMM群を用いて入力音声を認識することによって、高い認識率を得ることができるのである。

【0073】尚、上記実施の形態においては、上記話者クラスタ選択部3によって最適な話者クラスタを一つ選択するようにしているが、最適な話者クラスタを含む上位複数の話者クラスタを選択するようにしてもよい。例えば、尤度の上位からk個の話者クラスタを選択するとする。そうすると、切換部5によって切り換え選択されたk個の音響モデルの夫々に関して、尤度演算部4によって尤度演算が行われて、照合部8にk個の尤度系列が送られることになる。したがって、照合部8では、夫々の尤度系列に関して照合処理が行われ、最も大きい尤度を呈する単語/単語列が認識結果となるのである。

【0074】また、音声認識装置のハードウェア規模が大きく、計算量が許すのであれば、話者クラスタ選択部

3による話者クラスタ選択を行わず、尤度演算部において総ての話者クラスタの音響モデルを用いて尤度演算処理を実行するようにしてもよい。この場合、各音響モデルを適用して得られた尤度が最大値を呈する単語/単語列が認識結果となる。

【0075】<第2実施の形態>図2は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者正規化方式を用いた音声認識装置であり、話者正規化部26を有している。音声入力部21、音響分析部22、尤度演算部24、照合部28、言語モデル格納部29および出力部30は、図1に示す上記第1実施の形態における音声入力部1、音響分析部2、尤度演算部4、照合部8、辞書格納部9および出力部10と同様である。

【0076】上記話者正規化部26は、話者特徴抽出部25と周波数ワープ部23とから構成される。話者特徴抽出部25は、図1に示す上記第1実施の形態における話者特徴抽出部11と同様であり、入力された音響パラメータ系列に対して最大値を与えるGMMのクラスタ情報をGMM格納部(図示せず)から抽出して話者特徴とする。そして、得られたクラスタ情報から周波数伸縮係数 α を得、周波数ワープ部23に送出する。

【0077】そうすると、上記周波数ワープ部23は、この周波数伸縮係数 α を係数とする線形周波数ワーピング関数を用いて、入力音声の音声パラメータ系列を周波数ワープ(話者正規化)し、周波数ワープ後の音響パラメータ系列を尤度演算部24に送出するのである。そして、上記尤度演算部24では、周波数ワープされた音響パラメータ系列に対して、正規化不特定話者音響モデル格納部27に格納された不特定話者モデル(HMM)を作用させて、各音韻の状態毎に尤度を算出するのである。

【0078】ここで、上記正規化不特定話者音響モデル格納部27に格納される不特定話者モデルは、総ての学習話者を周波数伸縮によって $\alpha=1$ となるように正規化してから通常のHMM学習を行って作成される。尚、学習話者が多量に存在する場合には、全学習話者のうち、 $\alpha=1$ の話者およびその周辺の話者を正規化して学習の対象としてもよい。

【0079】上記第1実施の形態におけるGMM群作成時における話者クラスタリングの場合と同様に、音声認識時における話者正規化部26による話者正規化と、正規化不特定話者音響モデル格納部27に格納される不特定話者モデルの学習との場合にも、表1に示す推定/正規化処理の対象となる音素の分類に従って、以下のよう

に正規化対象とする音声区間を限定してもよい。

・音声認識時における話者正規化時…[c],[d],[e],[b])

・不特定話者モデルの学習時…[c],[d],[e],[b])

【0080】上述のように、本実施の形態においては、上記話者正規化部26によって入力話者を正規化する際

に、話者特徴抽出部25によって、上記第1実施の形態における話者特徴抽出部11の場合と同様にして、GMM格納部に各話者クラスタ毎に格納されたGMMを入力音響パラメータ系列に作用させて、最大尤度を与えるGMMのインデックス(係数 α)をクラスタ情報として求める。そして、周波数ワープ部23によって、上記クラスタ情報(係数 α)を係数とする線形周波数ワーピング関数を用いて入力音声の音響パラメータ系列を周波数ワープすることによって、話者正規化するようにしている。

【0081】その場合、上記話者特徴抽出部25が用いるGMM格納部には、上記第1実施の形態におけるGMM格納部6の場合と同様に、各学習話者の音声データを声道長伸縮係数 α に基づいてクラスタリングし、あるクラスタCの近傍のクラスタDに属する話者の音声データに基づいて上記クラスタCに属する音声データを生成し、各クラスタのGMMを用いて学習話者をクラスタリングし直したものが格納されている。したがって、話者特徴抽出部25は、少ない学習データから、発話内容に因らずに話者の特徴を良く表わすクラスタ情報を得ることができる。その結果、高い認識率を得ることができるのである。

【0082】<第3実施の形態>図3は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者適応方式を用いた音声認識装置であり、話者適応部36を有している。音声入力部31、音響分析部32、尤度演算部37、照合部39、言語モデル格納部40および出力部41は、図1に示す上記第1実施の形態における音声入力部1、音響分析部2、尤度演算部4、照合部8、言語モデル格納部9および出力部10と同様である。また、話者特徴抽出部33は、図2に示す上記第2実施の形態における話者特徴抽出部25と同様である。

【0083】上記話者適応部36は、上記話者特徴抽出部33と周波数ワープ部34とから構成される。話者特徴抽出部33は、上記第2実施の形態の場合と同様にして、入力された音響パラメータ系列に対して最大値を与えるGMMのクラスタ情報をGMM格納部(図示せず)から抽出して話者特徴とする。そして、得られたクラスタ情報から周波数伸縮係数 α を得、周波数ワープ部34に送出する。

【0084】そうすると、上記周波数ワープ部34は、この得られた周波数伸縮係数 α の逆数を係数とする線形周波数ワーピング関数を用いて、正規化不特定話者音響モデル格納部35に格納された不特定話者モデルを周波数ワープする。その場合の周波数ワープに際しては、上記第1実施の形態におけるGMM群作成時における話者クラスタリングの場合と同様に、表1に示す正規化処理の対象となる音素の分類に従って、以下のよう

に適応化対象とする音声区間を限定するのである。

・音声認識時における周波数ワープ時…[b],[c],

[d],[e]

但し、声道長の影響を受け難い[b]は変換しなくてもよい。

【0085】こうして周波数ワープされた不特定話者音響モデルは、話者適応モデル(HMM)として話者適応音響モデル格納部38に格納される。そうすると、尤度演算部37は、音響分析部32からの入力音声の音響パラメータ系列に対して、話者適応音響モデル格納部38に格納された話者適応モデルを作用させて、上述した尤度演算処理を行なうのである。

【0086】その場合、上記話者特徴抽出部33が用いるGMM格納部には、上記第1実施の形態におけるGM格納部6の場合と同様に、各学習話者の音声データを声道長伸縮係数 α に基づいてクラスタリングし、あるクラスタCの近傍のクラスタDに属する話者の音声データに基づいて上記クラスタCに属する音声データを生成し、各クラスタのGMMを用いて学習話者をクラスタリングし直したものが格納されている。したがって、話者特徴抽出部33は、少ない学習データから、発話内容に因らずに話者の特徴を良く表わすクラスタ情報を得ることができる。その結果、高い認識率を得ることができるのである。

【0087】尚、本実施の形態における上記話者適応音響モデル格納部38に格納する話者適応モデルの与え方には、上述の与え方の以外に、話者クラスタを用いる方法を採用してもよい。そして、この二通りの与え方を、音声認識装置の規模や入力音声データの量や質に応じて使い分けるのである。ここで、音声データの質とは尤度の上昇具合であり、話者特徴抽出部33は、上記二通りの与え方による尤度の上昇具合を見計らって、上昇の大きい方法を採用するのである。長いエンロール期間が許容できる音声認識装置の場合には、このような推定処理も可能となる。尚、上記話者クラスタを用いる方法においては、教師語彙を与える上記第1実施の形態における選択法[b]に基づいて話者クラスタを選択する。そして、選択された話者クラスタの音響モデルを話者適応モデルとして話者適応音響モデル格納部38に格納するのである。

【0088】尚、上述した各実施の形態においては、各学習話者の音声データを声道長伸縮係数 α に基づいてクラスタリングし、各クラスタのGMMを用いて学習話者をクラスタリングし直したGMM格納部を搭載した音声認識装置、および、上記クラスタリングしたGMMで成る話者特徴を用いて話者正規化あるいは話者適応を行う音声認識装置について説明している。しかしながら、この発明は、上記クラスタリングされたGMMの何れかを話者特徴として抽出する話者特徴抽出装置にも適用されるものである。

【0089】ところで、上記第1実施の形態における話者クラスタ選択部3、GMM格納部6およびGMM学習

部12による上記話者特徴抽出装置としての機能は、プログラム記録媒体に記録された話者特徴抽出処理プログラムによって実現される。上記プログラム記録媒体は、ROMでなるプログラムメディアである。または、外部補助記憶装置に装着されて読み出されるプログラムメディアであってもよい。尚、何れの場合においても、上記プログラムメディアから話者特徴抽出処理プログラムを読み出すプログラム読み出し手段は、上記プログラムメディアに直接アクセスして読み出す構成を有していてもよいし、RAM(ランダム・アクセス・メモリ)に設けられたプログラム記憶エリア(図示せず)にダウンロードして、上記プログラム記憶エリアにアクセスして読み出す構成を有していてもよい。尚、上記プログラムメディアからRAMのプログラム記憶エリアにダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。

10

20

30

【0090】ここで、上記プログラムメディアとは、本体側と分離可能に構成され、磁気テープやカセットテープ等のテープ系、フロッピー(登録商標)ディスク、ハードディスク等の磁気ディスクやCD(コンパクトディスク)・ROM、MO(光磁気)ディスク、MD(ミニディスク)、DVD(デジタルビデオディスク)等の光ディスクのディスク系、IC(集積回路)カードや光カード等のカード系、マスクROM、EPROM(紫外線消去型ROM)、EEPROM(電氣的消去型ROM)、フラッシュROM等の半導体メモリ系を含めた、固定的にプログラムを担持する媒体である。

【0091】また、上記各実施の形態における音声認識装置、音声合成装置および話者特徴抽出装置は、モデムを備えてインターネットを含む通信ネットワークと接続可能な構成を有していれば、上記プログラムメディアは、通信ネットワークからのダウンロード等によって流動的にプログラムを担持する媒体であっても差し支えない。尚、その場合における上記通信ネットワークからダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。または、別の記録媒体からインストールされるものとする。

【0092】尚、上記記録媒体に記録されるものはプログラムのみに限定されるものではなく、データも記録することが可能である。

40

【0093】

【発明の効果】以上より明らかなように、第1の発明の話者特徴抽出装置は、学習話者をクラスタリングするに際して、先ず、学習話者クラスタリング手段によって、各学習話者を標準話者に対する声道長の伸縮係数 α に基づいてクラスタリングし、音響モデル生成手段によって、各クラスタに属する話者集合毎にGMMを生成し、再クラスタリング手段によって、上記GMM群を用いて学習話者を再クラスタリングするので、各クラスタの初期値として声道長という大局的な特徴を明示的に与え

50

て、効率よくクラスタリングを行うことができる。さらに、GMMを用いて学習話者を再クラスタリングすることによって、発話内容に因らずに話者の特徴を良く表わすクラスタを得ることができ、声道長伸縮係数 α の抽出誤りも修復できる。

【0094】したがって、ループ学習手段によって、所定の条件を満たすまで上記GMMの生成と学習話者の再クラスタリングとを繰り返して得られたGMMが格納された音響モデル格納部を用いて、話者クラスタ選択部によって、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択することによって、発話内容に因らずに精度良く入力話者の特徴を抽出することができる。

【0095】また、1実施例の話者特徴抽出装置は、上記学習話者のクラスタリングに際して、音声サンプル豊富化手段によって、ある注目クラスタの隣接クラスタまたは近傍クラスタに属する話者の音声サンプルに対して周波数伸縮を行い、上記注目クラスタに属する音声サンプルを生成して編入するので、学習データ不足を補って、少ない発声データからでも各クラスタの音響モデルを精密に構築できる。したがって、少ない発声データでより精度良く入力話者の特徴を抽出できる。

【0096】また、1実施例の話者特徴抽出装置は、上記音声サンプル豊富化手段によって上記学習話者の音声サンプルに対して周波数伸縮を行う音声区間を、有音・無音の別および調音点に基づいて限定するので、声道長の差の影響を受け難い音素や無音部を上記周波数軸伸縮の対象外にして、声道長の差の影響を受け難い音素や無音部まで変形されることを防止できる。

【0097】また、1実施例の話者特徴抽出装置は、上記再クラスタリング手段によって上記学習話者を再クラスタリングする場合に、上記再クラスタリングの対象となる対象学習話者が再クラスタリングの前後に属しているクラスタの伸縮係数 α が所定値以上離れている場合には、当該対象学習話者を上記再クラスタリングの対象から外すので、声道長伸縮係数 α が極端に異なる話者同士が同じクラスタに属することを防止できる。

【0098】また、1実施例の話者特徴抽出装置は、サブクラスタ生成手段によって、上記ループ学習手段によって最終的にクラスタリングされた各クラスタに属する学習話者を更にクラスタリングしてサブクラスタを生成し、上記各サブクラスタに属する話者集合毎にGMMを生成し、上記音響モデル格納部は、上記サブクラスタ生成手段によって生成されたGMMの群を、各サブクラスタの伸縮係数 α に対応付けて格納するので、このサブクラスタによって声道長以外の要因にも対応することができ、より精密な話者特徴を抽出できる。

【0099】また、第2の発明の話者特徴抽出装置は、標準話者に対する声道長の伸縮係数 α に基づいて学習話者をクラスタリングし、各クラスタ毎のGMMの生成とそのGMM群を用いた学習話者の再クラスタリングとを

所定の条件を満たすまで繰り返し、最終的に生成されたGMM群を話者クラスタ毎に格納した音響モデル格納部を用いて、話者クラスタ選択部によって、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択して入力話者の特徴とするので、発話内容に拘らず精度良く入力話者の特徴を抽出することができる。

【0100】また、第3の発明の音声認識装置は、上記第1の発明あるいは第2の発明の話者特徴抽出装置の音響モデル格納部における話者クラスタ毎に生成されたHMMの群が格納されたHMM格納部から、切換部によって、上記話者特徴抽出装置で選択された話者クラスタのHMMを音声認識用の音響モデルとして切り換え選出するので、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わすHMMを用いて、入力話者の音声を正確に認識することができる。

【0101】また、第4の発明の音声認識装置は、上記第1の発明あるいは第2の発明に係る音響モデル格納部における話者クラスタ毎に生成されたHMMの群が格納されたHMM格納部を用いて、切換部によって、最大の尤度を呈するHMMを音声認識用の音響モデルとして切り換え選出するので、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わすHMMを用いて、入力話者の音声を正確に認識することができる。

【0102】また、第5の発明の音声認識装置は、話者正規化手段を、上記第1の発明あるいは第2の発明の話者特徴抽出装置と、入力話者の音声サンプルに基づいて上記話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α を用いて、上記入力音声のスペクトルの周波数軸を伸縮する周波数ワープ手段で構成したので、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わす声道長の伸縮係数 α を用いて、より標準話者の音響特徴量に近づくように話者正規化を行うことができる。したがって、高い音声認識率を得ることができる。

【0103】また、第6の発明の音声認識装置は、話者適応手段を、上記第1の発明あるいは第2の発明の話者特徴抽出装置と、入力話者の音声サンプルに基づいて上記話者特徴抽出装置によって選択されたGMMに対応する声道長の伸縮係数 α の逆数を用いて、音響モデルのスペクトルの周波数軸を伸縮する周波数ワープ手段で構成したので、発話内容に拘らずに少ない音声データで、上記入力話者の特徴を精度良く表わす声道長の伸縮係数 α の逆数を用いて、より入力話者の音響特徴量に近づくように話者適応を行うことができる。したがって、高い音声認識率を得ることができる。

【0104】また、第7の発明の話者特徴抽出方法は、各学習話者を標準話者に対する声道長の伸縮係数 α に基づいてクラスタリングし、各クラスタに属する話者集合毎にGMMを生成し、上記GMM群を用いて学習話者を

再クラスタリングするので、各クラスタの初期値として声道長という大局的な特徴を明示的に与えて、効率よくクラスタリングを行うことができる。さらに、GMMを用いて学習話者を再クラスタリングすることによって、発話内容に因らずに話者の特徴を良く表わすクラスタを得ることができ、声道長伸縮係数 α の抽出誤りも修復できる。

【0105】したがって、所定の条件を満たすまで上記GMMの生成と学習話者の再クラスタリングとを繰り返して得られたGMMを格納した音響モデル格納部を用いて、入力話者の音声サンプルに対して最大尤度を呈するGMMを選択することによって、発話内容に因らずに精度良く入力話者の特徴を抽出することができる。

【0106】また、第8の発明のプログラム記録媒体は、コンピュータを、上記第1の発明に係る学習話者クラスタリング手段、音響モデル生成手段、再クラスタリング手段、ループ学習手段、音響モデル格納部および話者クラスタ選択部として機能させる話者特徴抽出処理プログラムを記録しているので、上記第1の発明の場合と同様に、発話内容に因らずに精度良く入力話者の特徴を抽出することができる。

【図面の簡単な説明】

【図1】 この発明の話者特徴抽出装置を搭載したクラスタリング方式による音声認識装置におけるブロック図である。

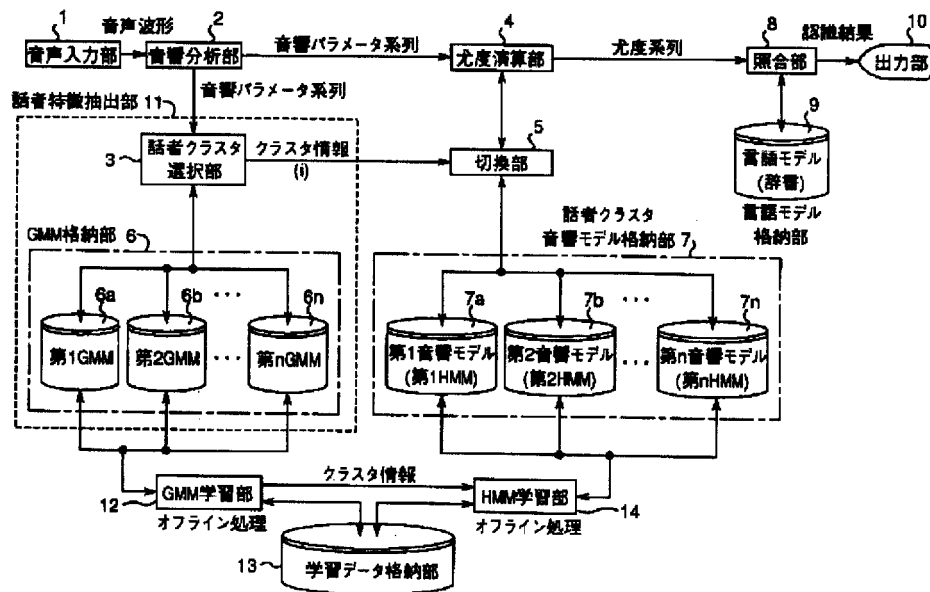
【図2】 図1とは異なる話者正規化方式による音声認識装置におけるブロック図である。

【図3】 図1および図2とは異なる話者適応方式による音声認識装置におけるブロック図である。

【符号の説明】

- 1, 2 1, 3 1…音声入力部、
- 2, 2 2, 3 2…音響分析部、
- 3…話者クラスタ選択部、
- 4, 2 4, 3 7…尤度(音韻類似度)演算部、
- 5…切換部、
- 6…GMM格納部、
- 7…話者クラスタ音響モデル格納部、
- 8, 2 8, 3 9…照合部、
- 9, 2 9, 4 0…言語モデル格納部、
- 10, 3 0, 4 1…出力部、
- 11, 2 5, 3 3…話者特徴抽出部、
- 12…GMM学習部、
- 13…学習データ格納部、
- 14…HMM学習部、
- 20 2 3, 3 4…周波数ワープ部、
- 26…話者正規化部、
- 27, 3 5…正規化不特定話者音響モデル格納部、
- 36…話者適応部、
- 38…話者適応音響モデル格納部。

【図1】



[illegible]

```

graph LR
    31[音声入力部] --> 32[音響分析部]
    32 -- "音響パラメータ系列" --> 37[尤度演算部]
    32 --> 33[話者特徴抽出部]
    33 -- "クラスタ情報" --> 34[周波数ワープ部]
    33 -- "係数 α" --> 34
    34 -- "変換" --> 38[(話者適応モデル(HMM))]
    35[(不正定話者モデル(HMM))] -- "正規化不正定話者音響モデル格納部" --> 34
    34 -- "周波数" --> 36[話者適応部]
    36 -- "話者適応パラメータ" --> 37
    37 -- "尤度系列" --> 39[照合部]
    39 -- "認識結果" --> 41[出力部]
    39 --> 40[(言語モデル辞書)]
    40 -- "言語モデル格納部" --> 39
  
```

Figure 1 is a block diagram of a speaker identification system. The system consists of the following components and data flow:

- 31 音声入力部 (Audio Input Unit):** Receives the audio input.
- 32 音響分析部 (Acoustic Analysis Unit):** Processes the input audio to generate an **音響パラメータ系列 (Acoustic Parameter Series)**, which is sent to the likelihood calculation unit (37).
- 33 話者特徴抽出部 (Speaker Feature Extraction Unit):** Extracts features from the audio analysis unit (32) to produce **クラスタ情報 (Cluster Information)** and a **係数 α (Coefficient α)**.
- 34 周波数ワープ部 (Frequency Warping Unit):** Uses the cluster information and coefficient α to perform a **変換 (Transformation)** on the **不正定話者モデル(HMM) (Indefinite Speaker Model (HMM))** (35) to generate a **話者適応モデル(HMM) (Speaker Adaptation Model (HMM))** (38).
- 35 不正定話者モデル(HMM) (Indefinite Speaker Model (HMM))** and **36 話者適応部 (Speaker Adaptation Unit):** The adaptation unit (36) uses the frequency warping unit (34) to produce a **話者適応パラメータ (Speaker Adaptation Parameter)** (37).
- 37 尤度演算部 (Likelihood Calculation Unit):** Combines the acoustic parameter series (32) and the speaker adaptation parameter (37) to produce an **尤度系列 (Likelihood Series)** (39).
- 39 照合部 (Verification Unit):** Compares the likelihood series (39) with a **言語モデル辞書 (Language Model Dictionary)** (40) to produce a **認識結果 (Recognition Result)** (41).
- 40 言語モデル辞書 (Language Model Dictionary):** A database used for verification.
- 41 出力部 (Output Unit):** Outputs the final recognition result.

5 3 5 C